

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号  
特許第7362929号  
(P7362929)

(45)発行日 令和5年10月17日(2023.10.17)

(24)登録日 令和5年10月6日(2023.10.6)

(51)国際特許分類 F I  
G 1 0 L 13/10 (2013.01) G 1 0 L 13/10 1 1 3 Z  
G 1 0 L 25/30 (2013.01) G 1 0 L 25/30

請求項の数 28 (全27頁)

(21)出願番号	特願2022-534694(P2022-534694)	(73)特許権者	502208397
(86)(22)出願日	令和1年12月10日(2019.12.10)		グーグル エルエルシー
(65)公表番号	特表2023-505670(P2023-505670 A)		Google LLC
(43)公表日	令和5年2月10日(2023.2.10)		アメリカ合衆国 カリフォルニア州 9 4 0 4 3 マウンテン ビュー アンフィシ
(86)国際出願番号	PCT/US2019/065566		アター パークウェイ 1 6 0 0
(87)国際公開番号	WO2021/118543		1 6 0 0 Amphitheatre P
(87)国際公開日	令和3年6月17日(2021.6.17)		arkway 9 4 0 4 3 Mounta
審査請求日	令和4年7月20日(2022.7.20)		in View, CA U . S . A .
		(74)代理人	100142907
			弁理士 本田 淳
		(72)発明者	クラーク、ロバート
			アメリカ合衆国 9 4 0 4 3 カリフォル
			ニア州 マウンテン ビュー アンフィシ
			アター パークウェイ 1 6 0 0
			最終頁に続く

(54)【発明の名称】 アテンションベースのクロックワーク階層型変分エンコーダ

(57)【特許請求の範囲】

【請求項1】

方法(400)であって、

データ処理ハードウェア(112)において、少なくとも1つの単語(240)を有するテキスト発話(310)を受信することであって、各単語(240)は少なくとも1つの音節(230)を有し、各音節(230)は少なくとも1つの音素(220)を有する、テキスト発話(310)を受信すること、

前記データ処理ハードウェア(112)によって、前記テキスト発話(310)のための発話埋め込み(204)であって、意図された韻律を表す前記発話埋め込み(204)を選択すること、

選択された発話埋め込み(204)を用いて、各音節(230)について、

前記データ処理ハードウェア(112)によって、前記音節(230)の各音素(220)の言語的特徴(222)に対するアテンション機構(340)によるアテンションに基づいて、前記音節(230)の韻律音節埋め込み(232, 234)を復号化することにより、前記音節(230)の持続時間(238)を予測すること、

前記データ処理ハードウェア(112)によって、前記音節(230)の予測された持続時間(238)に基づいて、複数の固定長予測フレーム(260)を生成することを含む、方法(400)。

【請求項2】

前記データ処理ハードウェア(112)によって、前記音節(230)の前記予測され

た持続時間(238)に基づいて前記音節(230)のピッチ輪郭(F0)を予測することをさらに含み、

前記複数の固定長予測フレーム(260)は、固定長予測ピッチフレーム(260, 260F0)を含み、各固定長予測ピッチフレーム(260F0)は、前記音節(230)の予測されたピッチ輪郭(F0)の一部を表す、請求項1に記載の方法(400)。

【請求項3】

前記選択された発話埋め込み(204)を用いて、各音節(230)について、

前記データ処理ハードウェア(112)によって、前記音節(230)の前記予測された持続時間(238)に基づいて、各音節(230)のエネルギー輪郭(C0)を予測すること、

前記データ処理ハードウェア(112)によって、対応する音節(230)の前記予測された持続時間(238)に基づいて、複数の固定長予測エネルギーフレーム(260, 260C0)を生成すること

をさらに含み、各固定長エネルギーフレーム(260C0)は、前記対応する音節(230)の予測されたエネルギー輪郭(C0)を表す、請求項1または2に記載の方法(400)。

【請求項4】

前記複数の固定長予測フレーム(260)は、前記音節(230)の固定長予測スペクトルフレーム(260, 260M0)を含む、請求項1~3のいずれか一項に記載の方法(400)。

【請求項5】

前記テキスト発話(310)の階層的言語構造(200)を表すネットワークは、

前記テキスト発話(310)の各単語(240)を含む第1レベルと、

前記テキスト発話(310)の各音節(230)を含む第2レベルと、

前記テキスト発話(310)の各音節(230)の各固定長予測フレーム(260)を含む第3レベルと

を含む、請求項1~4のいずれか一項に記載の方法(400)。

【請求項6】

前記階層的言語構造(200)を表す前記ネットワークの前記第1レベルは、前記テキスト発話(310)の各単語(240)を表す長短期記憶(LSTM)処理ブロックを含み、

前記階層的言語構造(200)を表す前記ネットワークの前記第2レベルは、前記テキスト発話(310)の各音節(230)を表すLSTM処理ブロックを含み、前記第2レベルの前記LSTM処理ブロックは、前記第1レベルのLSTM処理ブロックに対して相対的に、かつ前記第1レベルのLSTM処理ブロックよりも高速にクロックし、

前記階層的言語構造(200)を表す前記ネットワークの前記第3レベルは、各固定長予測フレーム(260)を表すLSTM処理ブロックを含み、前記第3レベルの前記LSTM処理ブロックは、前記第2レベルのLSTM処理ブロックに対して相対的に、かつ前記第2レベルのLSTM処理ブロックよりも高速にクロックする、請求項5に記載の方法(400)。

【請求項7】

前記音節(230)の前記持続時間(238)を予測することは、

前記音節(230)に関連付けられた各音素(220)について、

対応する音素(220)の1つまたは複数の言語的特徴(222)を符号化すること、

符号化された1つまたは複数の言語的特徴(222)を前記アテンション機構(340)に入力すること、

前記アテンション機構(340)の前記アテンションを前記韻律音節埋め込み(232, 234)に適用すること

を含む、請求項1~6のいずれか一項に記載の方法(400)。

【請求項8】

10

20

30

40

50

前記韻律音節埋め込み(232, 234)は、前記発話埋め込み(204)に対応するフレーム(210)に基づく第1の音節埋め込み(232)と、前記発話埋め込み(204)の1つまたは複数の音素(220)に関連付けられた音素言語的特徴(222)に基づく第2の音節埋め込み(234)とを含む、請求項1~7のいずれか一項に記載の方法(400)。

【請求項9】

前記データ処理ハードウェア(112)によって、複数の基準オーディオ信号(202)を含むトレーニングデータを受信することであって、各基準オーディオ信号(202)は、人間音声の音声発話を含み、かつ対応する韻律を有する、トレーニングデータを受信すること、

10

前記データ処理ハードウェア(112)によって、各基準オーディオ信号(202)を、前記基準オーディオ信号(202)の前記対応する韻律を表す対応する固定長発話埋め込み(204)に符号化することにより、韻律モデル(300)のためのディープニューラルネットワーク(200)をトレーニングすること

をさらに含む、請求項1~8のいずれか一項に記載の方法(400)。

【請求項10】

前記データ処理ハードウェア(112)によって、フレームベースの音節埋め込み(232)および単音特徴ベースの音節埋め込み(234)で、複数の言語ユニット(220, 230, 240, 250)の言語的特徴(222, 236, 242, 252)を符号化することにより、前記選択された発話埋め込み(204)を生成することをさらに含む、請求項9に記載の方法(400)。

20

【請求項11】

前記発話埋め込み(204)は、固定長の数値ベクトルを含む、請求項1~10のいずれか一項に記載の方法(400)。

【請求項12】

前記アテンション機構(340)の前記アテンションは、位置ベースのアテンションを含む、請求項1~11のいずれか一項に記載の方法(400)。

【請求項13】

前記位置ベースのアテンションは、単調にシフトする、位置に敏感なアテンションを含み、前記単調にシフトする、位置に敏感なアテンションは、それぞれの音節(230)の音素情報のウィンドウによって定義される、請求項12に記載の方法(400)。

30

【請求項14】

前記アテンション機構(340)は、トランスフォーマーを含む、請求項1~11のいずれか一項に記載の方法(400)。

【請求項15】

データ処理ハードウェア(112)と、  
前記データ処理ハードウェア(112)と通信するメモリハードウェア(114)とを備えるシステム(100)であって、前記メモリハードウェア(114)は、前記データ処理ハードウェア(112)上で実行されると前記データ処理ハードウェア(112)に動作を実行させる命令を格納しており、前記動作は、

40

少なくとも1つの単語(240)を有するテキスト発話(310)を受信することであって、各単語(240)は少なくとも1つの音節(230)を有し、各音節(230)は少なくとも1つの音素(220)を有する、テキスト発話(310)を受信すること、

前記テキスト発話(310)のための発話埋め込み(204)であって、意図された韻律を表す前記発話埋め込み(204)を選択すること、

選択された発話埋め込み(204)を用いて、各音節(230)について、

前記音節(230)の各音素(220)の言語的特徴(222)に対するアテンション機構(340)によるアテンションに基づいて、前記音節(230)の韻律音節埋め込み(232, 234)を復号化することにより、前記音節(230)の持続時間(238)を予測すること、

50

前記音節（２３０）の予測された持続時間（２３８）に基づいて、複数の固定長予測フレーム（２６０）を生成すること  
を含む、システム（１００）。

【請求項１６】

前記動作は、

前記音節（２３０）の前記予測された持続時間（２３８）に基づいて前記音節（２３０）のピッチ輪郭（Ｆ０）を予測すること

をさらに含み、

前記複数の固定長予測フレーム（２６０）は、固定長予測ピッチフレーム（２６０，２６０Ｆ０）を含み、各固定長予測ピッチフレーム（２６０Ｆ０）は、前記音節（２３０）の予測されたピッチ輪郭（Ｆ０）の一部を表す、請求項１５に記載のシステム（１００）。

10

【請求項１７】

前記動作は、前記選択された発話埋め込み（２０４）を用いて、各音節（２３０）について、

前記音節（２３０）の前記予測された持続時間（２３８）に基づいて、各音節（２３０）のエネルギー輪郭（Ｃ０）を予測すること、

対応する音節（２３０）の前記予測された持続時間（２３８）に基づいて、複数の固定長予測エネルギーフレーム（２６０，２６０Ｃ０）を生成すること

をさらに含み、各固定長エネルギーフレーム（２６０Ｃ０）は、前記対応する音節（２３０）の予測されたエネルギー輪郭（Ｃ０）を表す、請求項１５または１６に記載のシステム（１００）。

20

【請求項１８】

前記複数の固定長予測フレーム（２６０）は、前記音節（２３０）の固定長予測スペクトルフレーム（２６０，２６０Ｍ０）を含む、請求項１５～１７のいずれか一項に記載のシステム（１００）。

【請求項１９】

前記テキスト発話（３１０）の階層的言語構造（２００）を表すネットワークは、

前記テキスト発話（３１０）の各単語（２４０）を含む第１レベルと、

前記テキスト発話（３１０）の各音節（２３０）を含む第２レベルと、

前記テキスト発話（３１０）の各音節（２３０）の各固定長予測フレーム（２６０）を含む第３レベルと

30

を含む、請求項１５～１８のいずれか一項に記載のシステム（１００）。

【請求項２０】

前記階層的言語構造（２００）を表す前記ネットワークの前記第１レベルは、前記テキスト発話（３１０）の各単語（２４０）を表す長短期記憶（ＬＳＴＭ）処理ブロックを含み、

前記階層的言語構造（２００）を表す前記ネットワークの前記第２レベルは、前記テキスト発話（３１０）の各音節（２３０）を表すＬＳＴＭ処理ブロックを含み、前記第２レベルの前記ＬＳＴＭ処理ブロックは、前記第１レベルのＬＳＴＭ処理ブロックに対して相対的に、かつ前記第１レベルのＬＳＴＭ処理ブロックよりも高速にクロックし、

40

前記階層的言語構造（２００）を表す前記ネットワークの前記第３レベルは、各固定長予測フレーム（２６０）を表すＬＳＴＭ処理ブロックを含み、前記第３レベルの前記ＬＳＴＭ処理ブロックは、前記第２レベルのＬＳＴＭ処理ブロックに対して相対的に、かつ前記第２レベルのＬＳＴＭ処理ブロックよりも高速にクロックする、請求項１９に記載のシステム（１００）。

【請求項２１】

前記音節（２３０）の前記持続時間（２３８）を予測することは、

前記音節（２３０）に関連付けられた各音素（２２０）について、

対応する音素（２２０）の１つまたは複数の言語的特徴（２２２）を符号化すること、

符号化された１つまたは複数の言語的特徴（２２２）を前記アテンション機構（３４

50

0)に入力すること、

前記アテンション機構(340)の前記アテンションを前記韻律音節埋め込み(232, 234)に適用すること

を含む、請求項15~20のいずれか一項に記載のシステム(100)。

【請求項22】

前記韻律音節埋め込み(232, 234)は、前記発話埋め込み(204)に対応するフレーム(210)に基づく第1の音節埋め込み(232)と、前記発話埋め込み(204)の1つまたは複数の音素(220)に関連付けられた音素言語的特徴(222)に基づく第2の音節埋め込み(234)とを含む、請求項15~21のいずれか一項に記載のシステム(100)。

10

【請求項23】

前記動作は、

複数の基準オーディオ信号(202)を含むトレーニングデータを受信することであって、各基準オーディオ信号(202)は、人間音声の音声発話を含み、かつ対応する韻律を有する、トレーニングデータを受信すること、

各基準オーディオ信号(202)を、前記基準オーディオ信号(202)の前記対応する韻律を表す対応する固定長発話埋め込み(204)に符号化することにより、韻律モデル(300)のためのディープニューラルネットワーク(200)をトレーニングすることをさらに含む、請求項15~22のいずれか一項に記載のシステム(100)。

【請求項24】

20

前記動作は、フレームベースの音節埋め込み(232)および単音特徴ベースの音節埋め込み(234)で、複数の言語ユニット(220, 230, 240, 250)の言語的特徴(222, 236, 242, 252)を符号化することにより、前記選択された発話埋め込み(204)を生成することをさらに含む、請求項23に記載のシステム(100)。

【請求項25】

前記発話埋め込み(204)は、固定長の数値ベクトルを含む、請求項15~24のいずれか一項に記載のシステム(100)。

【請求項26】

前記アテンション機構(340)の前記アテンションは、位置ベースのアテンションを含む、請求項15~25のいずれか一項に記載のシステム(100)。

30

【請求項27】

前記位置ベースのアテンションは、単調にシフトする、位置に敏感なアテンションを含み、前記単調にシフトする、位置に敏感なアテンションは、それぞれの音節(230)の音素情報のウィンドウによって定義される、請求項26に記載のシステム(100)。

【請求項28】

前記アテンション機構(340)は、トランスフォーマーを含む、請求項15~25のいずれか一項に記載のシステム(100)。

【発明の詳細な説明】

【技術分野】

40

【0001】

本開示は、アテンションベースのクロックワーク階層型変分エンコーダに関する。

【背景技術】

【0002】

音声合成システムは、テキスト音声変換(text-to-speech, TTS)モデルを使用して、テキスト入力から音声を生産する。生成/合成された音声は、意図された韻律(表現力)を備えた人間の音声のように聞こえながら(自然さ)、メッセージを正確に伝える(分かりやすさ)必要がある。従来の波形接続型(concatenative)およびパラメトリック合成モデルは分かりやすい音声を提供でき、音声のニューラルモデリングの最近の進歩により、合成音声の自然さが大幅に向上したが、既存のTTSモ

50

デルのほとんどは韻律のモデリングに効果がないため、重要なアプリケーションで使用される合成音声に表現力が不足する原因となっている。例えば、会話アシスタントおよび長文リーダー (long-form reader) などのアプリケーションでは、イントネーション、強勢、リズムおよびスタイルのような、テキスト入力では伝達されない韻律的な特徴を入力することで、リアルな音声を生成することが望ましい。例えば、単純なステートメントは、ステートメントが質問であるか、質問への回答であるか、ステートメントに不確実性があるか、または入力テキストによって指定されていない環境またはコンテキストに関するその他の意味を伝えるかどうかに応じて、さまざまな方法で話すことができる。

#### 【発明の概要】

10

#### 【0003】

本開示の一態様は、アテンション (attention) ベースのクロックワーク階層型変分エンコーダ (clockwork hierarchical variational encoder) のための方法を提供する。方法は、データ処理ハードウェアにおいて、少なくとも1つの単語を有するテキスト発話を受信することを含み、各単語は少なくとも1つの音節を有し、各音節は少なくとも1つの音素を有している。方法は、データ処理ハードウェアによって、テキスト発話のための発話埋め込みを選択することも含む。発話埋め込みは、意図された韻律を表す。方法は、選択された発話埋め込みを用いて、各音節について、データ処理ハードウェアによって、音節の各音素の言語的特徴に対するアテンション機構 (attention mechanism) によるアテンションに基づいて、音節の韻律音節埋め込みを復号化することにより、音節の持続時間 (duration) を予測することを含む。方法は、選択された発話埋め込みを用いて、各音節について、データ処理ハードウェアによって、音節の予測された持続時間に基づいて、複数の固定長予測フレームを生成することを含む。

20

#### 【0004】

本開示の実装形態は、以下の任意選択の特徴のうちの1つまたは複数を含んでよい。いくつかの実装形態では、方法は、データ処理ハードウェアによって、音節の予測された持続時間に基づいて音節のピッチ輪郭を予測することを含む。この実装形態において、複数の固定長予測フレームが固定長予測ピッチフレームを含む場合、各固定長予測ピッチフレームは、音節の予測されたピッチ輪郭の一部を表している。

30

#### 【0005】

いくつかの例において、方法は、選択された発話埋め込みを用いて、各音節について、データ処理ハードウェアによって、音節の予測された持続時間に基づいて、各音節のエネルギー輪郭を予測することを含む。この例では、方法は、データ処理ハードウェアによって、対応する音節の予測された持続時間に基づいて、複数の固定長予測エネルギーフレームを生成することも含む、各固定長エネルギーフレームは、対応する音節の予測されたエネルギー輪郭を表す。複数の固定長予測フレームは、音節の固定長予測スペクトルフレームを含み得る。

#### 【0006】

いくつかの構成では、テキスト発話の階層的言語構造を表すネットワークは、テキスト発話の各単語を含む第1レベルと、テキスト発話の各音節を含む第2レベルと、テキスト発話の各音節の各固定長予測フレームを含む第3レベルとを含む。ここで、階層的言語構造を表すネットワークの第1レベルは、テキスト発話の各単語を表す長短期記憶 (LSTM) 処理ブロックを含み得る。階層的言語構造を表すネットワークの第2レベルは、テキスト発話の各音節を表す LSTM 処理ブロックを含んでよく、第2レベルの LSTM 処理ブロックは、第1レベルの LSTM 処理ブロックに対して相対的に、かつ第1レベルの LSTM 処理ブロックよりも高速にクロックする。階層的言語構造を表すネットワークの第3レベルは、各固定長予測フレームを表す LSTM 処理ブロックを含んでよく、第3レベルの LSTM 処理ブロックは、第2レベルの LSTM 処理ブロックに対して相対的に、かつ第2レベルの LSTM 処理ブロックよりも高速にクロックする。

40

50

## 【0007】

いくつかの構成では、音節の長さを予測することは、音節に関連付けられた各音素について、対応する音素の1つまたは複数の言語的特徴を符号化すること、符号化された1つまたは複数の言語的特徴をアテンション機構に入力すること、アテンション機構のアテンションを韻律音節埋め込みに適用する ( *applying* ) ことを含む。韻律音節埋め込みは、発話埋め込みに対応するフレームに基づく第1の音節埋め込みと、発話埋め込みの1つまたは複数の音素に関連付けられた音素言語的特徴に基づく第2の音節埋め込みとを含み得る。

## 【0008】

いくつかの例では、方法は、データ処理ハードウェアによって、複数の基準オーディオ信号を含むトレーニングデータを受信すること、データ処理ハードウェアによって、各基準オーディオ信号の対応する韻律を表す対応する固定長発話埋め込みに符号化することにより、韻律モデルのためのディープニューラルネットワークをトレーニングすることを含む。この例では、各基準オーディオ信号は、人間音声の音声発話を含み、かつ対応する韻律を有している。ここで、方法は、データ処理ハードウェアによって、フレームベースの音節埋め込みおよび単音 ( *phone* ) 特徴ベースの音節埋め込みで、複数の言語ユニットの言語的特徴を符号化することにより、選択された発話埋め込みを生成することを含み得る。発話埋め込みは、固定長の数値ベクトルを含んでいてよい。

10

## 【0009】

いくつかの実装形態では、アテンション機構のアテンションは、位置ベースのアテンションを含む。位置ベースのアテンションは、単調にシフトする、位置に敏感な ( *location sensitive* ) アテンションを含み、単調にシフトする、位置に敏感なアテンションは、それぞれの音節の音素情報のウィンドウによって定義される。アテンション機構は、トランスフォーマー ( *transformer* ) を含み得る。

20

## 【0010】

本開示の別の態様は、アテンションベースのクロックワーク階層型変分エンコーダのためのシステムを提供する。システムは、データ処理ハードウェアおよびデータ処理ハードウェアと通信するメモリハードウェアを含む。メモリハードウェアは、データ処理ハードウェア上で実行されると、データ処理ハードウェアに動作を実行させる命令を格納している。動作は、少なくとも1つの単語を有するテキスト発話を受信することを含み、各単語は少なくとも1つの音節を有し、各音節は少なくとも1つの音素を有している。動作は、テキスト発話のための発話埋め込みを選択することも含み、発話埋め込みは、意図された韻律を表す。動作は、選択された発話埋め込みを用いて、各音節について、音節の各音素の言語的特徴に対するアテンション機構によるアテンションに基づいて、音節の韻律音節埋め込みを復号化することにより、音節の持続時間を予測することをさらに含む。動作は、選択された発話埋め込みを用いて、各音節について、音節の予測された持続時間に基づいて、複数の固定長予測フレームを生成することも含む。

30

## 【0011】

この態様は、以下の任意選択の特徴のうちの1つまたは複数を含んでいてよい。いくつかの構成では、動作は、音節の予測された持続時間に基づいて音節のピッチ輪郭を予測することを含み、複数の固定長予測フレームが固定長予測ピッチフレームを含む場合、各固定長予測ピッチフレームは、音節の予測されたピッチ輪郭の一部を表している。動作は、選択された発話埋め込みを用いて、各音節について、音節の予測された持続時間に基づいて、各音節のエネルギー輪郭を予測すること、対応する音節の予測された持続時間に基づいて、複数の固定長予測エネルギーフレームを生成することを含んでいてもよく、各固定長エネルギーフレームは、対応する音節の予測されたエネルギー輪郭を表す。複数の固定長予測フレームは、音節の固定長予測スペクトルフレームを含み得る。

40

## 【0012】

いくつかの例では、テキスト発話の階層的言語構造を表すネットワークは、テキスト発話の各単語を含む第1レベルと、テキスト発話の各音節を含む第2レベルと、テキスト発

50

話の各音節の各固定長予測フレームを含む第3レベルとを含む。ここで、階層的言語構造を表すネットワークの第1レベルは、テキスト発話の各単語を表す長短期記憶(LSTM)処理ブロックを含み得る。階層的言語構造を表すネットワークの第2レベルは、テキスト発話の各音節を表すLSTM処理ブロックを含んでいてよく、第2レベルのLSTM処理ブロックは、第1レベルのLSTM処理ブロックに対して相対的に、かつ第1レベルのLSTM処理ブロックよりも高速にクロックする。階層的言語構造を表すネットワークの第3レベルは、各固定長予測フレームを表すLSTM処理ブロックを含んでいてよく、第3レベルのLSTM処理ブロックは、第2レベルのLSTM処理ブロックに対して相対的に、かつ第2レベルのLSTM処理ブロックよりも高速にクロックする。

#### 【0013】

いくつかの実装形態では、音節の長さを予測することは、音節に関連付けられた各音素について、対応する音素の1つまたは複数の言語的特徴を符号化すること、符号化された1つまたは複数の言語的特徴をアテンション機構に入力すること、アテンション機構のアテンションを韻律音節埋め込みに適用することを含む。韻律音節埋め込みは、発話埋め込みに対応するフレームに基づく第1の音節埋め込みと、発話埋め込みの1つまたは複数の音素に関連付けられた音素言語的特徴に基づく第2の音節埋め込みとを含み得る。

#### 【0014】

いくつかの構成では、動作は、複数の基準オーディオ信号を含むトレーニングデータを受信することを含み、各基準オーディオ信号は、人間音声の音声発話を含み、かつ対応する韻律を有する。この構成では、動作は、各基準オーディオ信号を、基準オーディオ信号の対応する韻律を表す対応する固定長発話埋め込みに符号化することにより、韻律モデルのためのディープニューラルネットワークをトレーニングすることも含む。ここで、動作は、フレームベースの音節埋め込みおよび単音特徴ベースの音節埋め込みで、複数の言語ユニットの言語的特徴を符号化することにより、選択された発話埋め込みを生成することを含み得る。発話埋め込みは、固定長の数値ベクトルを含んでいてよい。

#### 【0015】

いくつかの例では、アテンション機構のアテンションは、位置ベースのアテンションを含む。ここで、位置ベースのアテンションは、単調にシフトする、位置に敏感なアテンションを含み、単調にシフトする、位置に敏感なアテンションは、それぞれの音節の音素情報のウィンドウによって定義される。アテンション機構は、トランスフォーマーを含み得る。

#### 【0016】

本開示の1つまたは複数の実施形態の詳細は、添付の図面および以下の説明に記載されている。他の態様、特徴、および利点は、説明および図面、ならびに特許請求の範囲から明らかになるであろう。

#### 【図面の簡単な説明】

#### 【0017】

【図1】図1は、テキスト発話の韻律表現を予測する際に使用する制御可能な韻律モデルを提供するためにディープニューラルネットワークをトレーニングするための例示的なシステムの概略図である。

【図2A】図2Aは、基準オーディオ信号の韻律を固定長の発話埋め込みに符号化するための階層的言語構造の概略図である。

【図2B】図2Bは、テキスト発話の韻律表現を予測するために発話埋め込みを使用する階層的言語構造の概略図である。

【図2C】図2Cは、テキスト発話の韻律表現を予測するために発話埋め込みを使用する階層的言語構造の概略図である。

【図3A】図3Aは、テキスト発話の音節特性を予測するための例示的なオートエンコーダの概略図である。

【図3B】図3Bは、テキスト発話の音節特性を予測するための例示的なオートエンコーダの概略図である。

10

20

30

40

50

【図 3 C】図 3 C は、テキスト発話の音節特性を予測するための例示的なオートエンコーダの概略図である。

【図 3 D】図 3 D は、符号化された音素状態を形成する単音レベルの特徴にアテンションを与えるように構成された例示的なアテンション機構の概略図である。

【図 4】図 4 は、受信されたテキスト発話の表現を予測する方法のための動作の例示的な配置のフローチャートである。

【図 5】図 5 は、本明細書で説明されるシステムおよび方法を実装するために使用され得る例示的なコンピューティングデバイスの概略図である。

【発明を実施するための形態】

【0018】

さまざまな図面の同様の参照記号は、同様の要素を示す。

音声合成システムでしばしば使用されるテキスト音声変換 (Text-to-Speech, TTS) モデルは、概して、ランタイムにおいて、いかなる基準音響表現もなしでテキスト入力のみを与えられ、本物らしく聞こえる合成音声を生成するために、テキスト入力によって提供されない多くの言語的ファクターを帰属させる必要がある。これらの言語的ファクターのサブセットは、まとめて韻律と呼ばれ、イントネーション (ピッチの変化)、強勢 (強勢のある音節対無強勢の音節)、音の持続時間、ラウドネス、トーン、リズム、および音声のスタイルを含んでよい。韻律は、音声の感情状態、音声の形式 (例えば、ステートメント、質問、コマンドなど)、音声の皮肉または嫌味の存在、音声の認識における不確実性、または入力テキストの文法または語彙選択によって符号化されることが不可能な他の言語要素を示し得る。したがって、大きな韻律変化に関連付けられた特定のテキスト入力は、ピッチおよび発話継続時間の局所的な変化を伴う合成音声を生成して、異なるセマンティック意味を伝えることができ、また、全体的なピッチ軌道のグローバルな変化を伴う合成音声を生成して、異なる気分および感情を伝えることもできる。

【0019】

ニューラルネットワークモデルは、テキスト入力では提供されない韻律に対応する言語的ファクターを予測することにより、音声を確実に合成する可能性を提供する。その結果、オーディオブックのナレーション、ニュースリーダー、ボイスデザインソフトウェア、および会話アシスタントなどの複数のアプリケーションが、単調には響かない本物らしく聞こえる合成音声を生成することができる。本明細書の実装形態は、音声発話に対応する基準オーディオ信号を音声発話の韻律を表す発話埋め込みに符号化するためのエンコーダ部分と、発話埋め込みを復号化して、音節の持続時間と、各音節のピッチおよびエネルギー輪郭とを予測するデコーダ部分とを有する変分オートエンコーダ (VAE) を含むニューラルネットワークモデルを対象としている。

【0020】

エンコーダ部分は、発話を表す言語的特徴を条件とする多数の基準オーディオ信号を符号化することによって、韻律を表す発話埋め込みをトレーニングすることができる。言語的特徴は、各音素の個々の音、各音節が強勢を有するかまたは無強勢であるか、発話における各単語のタイプ (例えば、名詞 / 形容詞 / 動詞) および / または単語の位置、ならびに発話が質問であるかまたはフレーズであるかを含み得るが、これらに限定されない。各発話埋め込みは、固定長の数値ベクトルによって表される。いくつかの実装形態では、固定長の数値ベクトルは、256 に等しい値を含む。しかしながら、他の実装形態では、256 より大きいまたは小さい値 (例えば、128) を有する固定長の数値ベクトルが使用されてもよい。デコーダ部分は、固定長発話埋め込みを、第 1 のデコーダを介して音節持続時間のシーケンスに復号化するとともに、音節持続時間を使用してピッチおよびエネルギーの固定長フレーム (例えば、5 ミリ秒) のシーケンスに復号化することができる。トレーニング中、デコーダ部分によって予測された音節持続時間と、ピッチおよびエネルギーの固定長フレームとは、固定長発話埋め込みに関連付けられた基準オーディオ信号からサンプリングされた音節持続時間と、ピッチおよびエネルギーの固定長フレームとに厳密に一致する。

10

20

30

40

50

## 【 0 0 2 1 】

本開示のVAEは、長短期記憶(LSTM)ブロックの階層的なスタックされた層を組み込んだクロックワーク階層型変分オートエンコーダ(Clockwork Hierarchical Variational Autoencoder, CHiVE)を含み、LSTMブロックの各層は、発話の構造を組み込んでいる。ここで、各LSTMブロックは、1つまたは複数のLSTMセルに分割することができる。発話は、音素、音節、単語、フレーズ、または文などの言語ユニットのいずれか1つまたは組み合わせに分割できるため、LSTMブロックは、そのようなユニットを表す1つまたは複数の層を含み得る。例えば、LSTMブロックは、音素、音節、単語、フレーズ、または文を表す1つまたは複数の層を含むLSTMセルを含む。さらに、LSTMセルのスタックされた層の階層は、階層的入力データの長さに合わせて可変的にクロックされる。例えば、入力データが、3音節の単語と、それに続く4音節の単語を含む場合、CHiVEの音節層は、第1の入力単語について、単語層の単一クロックに対し3回クロックし、次いで、音節層は、第2の単語について、単語層の後続の単一クロックに対しさらに4回クロックする。このように、所与のLSTMセルに関連付けられたメモリが約0.5秒(すなわち、5ミリ秒のフレームレートで100回ステップ)しか有効でなく、したがって、音声の2音節または3音節分のLSTMセルメモリしか提供できないフレームベースの技術を使うのではなく、CHiVEの音素、単語、および音節層は、それぞれ音素、単語、音節でクロックして、スタックされた層のLSTMセルに、過去100単語、音節または音素にわたるメモリを与える。追加的または代替的に、CHiVEは、発話構造を表すために階層的な層にLSTM構造を使用する代わりに、他の形式のニューラルネットワーク(NN)またはリカレントニューラルネットワーク(RNN)を使用するように適合されていてもよい。

10

20

## 【 0 0 2 2 】

推論中、CHiVEは、テキスト発話を受信し、テキスト発話のために発話埋め込みを選択するように構成されている。受信されたテキスト発話は、少なくとも1つの単語を有し、各単語は、少なくとも1つの音節を有し、各音節は、少なくとも1つの音素を有している。テキスト発話は、発話から合成音声を生成するための適切な韻律を導くためのコンテキスト、セマンティック情報、および語用論情報を欠いているため、CHiVEは、その選択された発話埋め込みを潜在変数として使用して、意図された韻律を表す。その後、CHiVEは、選択された発話埋め込みを使用して、音節に含まれる各音素の言語的特徴を、その音節に対する対応する韻律音節埋め込みで符号化することにより、各音節の持続時間を予測し、音節に対する予測された持続時間に基づいて各音節のピッチを予測する。最後に、CHiVEは、各固定長ピッチフレームが音節の予測ピッチを表すように、各音節の予測された持続時間に基づいて複数の固定長ピッチフレームを生成するように構成されている。CHiVEは、同様に、音節の予測された持続時間に基づいて各音節のエネルギー(例えば、ラウドネス)を予測し、各々が音節の予測されたエネルギーを表す複数の固定長エネルギーフレームを生成し得る。固定長ピッチおよび/またはエネルギーフレームは、TTSシステムのユニット選択モデルまたはウェーブネットモデルに提供されて、入力固定長発話埋め込みによって提供される意図された韻律を備えた合成音声を生成することができる。

30

40

## 【 0 0 2 3 】

一般に、いくつかの音声合成システムは、2つのフェーズに分けられる場合がある。発話に含まれる音声に影響を与えるファクターを識別する言語仕様を生成する第1フェーズ、および言語仕様を使用して合成音声の波形を生成する第2フェーズである。いくつかの例では、音声合成システムの別の側面が音声を生成するために使用する言語特性を予測する代わりに、CHiVEは、スペクトルフレーム(例えば、メルフレームなどの複数の固定長スペクトルフレーム)を予測するように代替的または追加的に構成されている。スペクトルフレームを予測することにより、CHiVEは、言語仕様から波形を生成するための音声合成システムのさらなる処理を最小化することができる。ここで、CHiVEは、選択された発話埋め込みを使用して、音節に含まれる各音素の言語的特徴を、その音節に

50

対する対応する韻律音節埋め込みで符号化することにより、各音節の持続時間を予測する。各音節に対する予測された持続時間を使用して、CHiVEは、複数の固定長スペクトルフレームを生成することができる。例えば、第1のデコーダは、各音節に対する予測された持続時間を生成し、これは、第2のデコーダが各音節に対して生成すべきスペクトルフレームの数を示す。次に、第2のデコーダは、各音節に対する予測された持続時間によって示されるスペクトルフレームの数をシーケンス復号化する。次に、スペクトルフレームは、音声合成システムによって使用されて、合成音声を生成することができる。例えば、スペクトルフレームは、ニューラルボコーダに提供される。

#### 【0024】

残念ながら、特定の言語ユニットに依存すると、ピッチ、エネルギー、またはスペクトルフレームのような言語ファクターを予測するのに支障が生じる可能性がある。特に、音声サンプル（例えば、音声合成システムをトレーニングする録音されたサンプル）は、音素シーケンスが不正確な可能性がある。これらの不正確さのために、音節アライメント（例えば、音節の境界）を使用し、次いで音素の境界を使用する方が、信頼性が高い場合がある。音素の場合、音声サンプルから予想される各音素を明確に定義するのは容易ではなく単純でもない場合がある。人は3つの音素を発音しようとし得るが、文脈やその人の方言によっては、各音素を明確または正確に発音できない場合がある。収集されたサンプルでは、サンプルを提供する人によって、または同じ人からのサンプルの間で、音素（音素レベル）の方言に違いがある場合がある。同じ人でも、文章またはその人のコミュニケーション速度に応じて単語を異なって発音する場合がある。素早く話す場合、人は音素をブレンドしたり、または音素を完全に脱落させたりすることがある。

#### 【0025】

別の例として、しばしば、地域の方言では、その地域内でよく使用される単語の発音が明確でない場合がある。説明のために、これは都市名で起こり得る。例えば、エジンバラ（Edinburgh）市には少なくとも3つの音節、すなわち「Ed-in-burgh」を有しているが、早口や地元の方言で表現すると、「Em-bra」と発音され、基本的には2音節に短縮される。これらの2つの音節では、エジンバラの発音に存在すると予想されるすべての音素を一致させることは困難である（例えば、3音節の形式）。これらの発音の違いは、短い基本的な単語または複数の音節を持つ大きな単語で発生する可能性がある。例えば、「th-ahh-t」の発音を有する「that」のような単純な単語は、「thuh-t」に短縮される場合がある。音声の発音が最適ではない場合（例えば、曖昧になったり急いだりした場合）、サンプルに存在する音素シーケンスは必ずしも正確ではない。モデルがこれらのタイプの不正確な音素サンプルに音素を割り当てようとした場合、モデルはエラーを起こす危険性がある。

#### 【0026】

信頼できないデータの原因となるこれらの音素の問題のいくつかを克服するために、CHiVEは、音素アライメントに依存しないアプローチを使用することができる。例えば、CHiVEは、音素の数に依存せず、音節のフレーム数を予測しようとする。ここで、CHiVEは、音素ごとの個々のフレーム数について何も知らない可能性がある。いくつかの例では、音節の音声コンテンツに関する情報は、音素の言語的特徴（例えば、音節に関する言語情報）に基づいている。これらの例では、CHiVEは、音節情報（例えば、音節あたりのフレーム数）を予測するために、符号化されてアテンション機構に提供される音素の言語的特徴を使用する。

#### 【0027】

図1は、制御可能な韻律モデル300を提供するためにディープニューラルネットワーク200をトレーニングし、韻律モデル300を使用してテキスト発話310の韻律表現302を予測するための例示的なシステム100を示す。システム100は、データ処理ハードウェア112、およびデータ処理ハードウェア112と通信し、データ処理ハードウェア112に動作を実行させる命令を格納しているメモリハードウェア114を有するコンピューティングシステム110を含む。いくつかの実装形態では、コンピューティン

10

20

30

40

50

グシステム 110 (例えば、データ処理ハードウェア 112) は、入力テキスト発話 310 からの合成音声 122 の韻律を制御するためにテキスト音声変換 (text-to-speech, TTS) システム 120 に、トレーニングされたディープニューラルネットワーク 200 に基づく韻律モデル 300 を提供する。入力テキスト発話 310 は、合成音声 122 の適切な韻律を導くためのコンテキスト、セマンティクス、および語用論を伝達する方法を有しないので、韻律モデル 300 は、テキスト発話 310 から抽出された言語的特徴でモデル 300 を条件付け、固定長発話埋め込み 204 を、テキスト発話 310 の意図する韻律を表す潜在変数として使用することによって、入力テキスト発話 310 の韻律表現 302 を予測することができる。いくつかの例では、コンピューティングシステム 110 は、TTS システム 120 を実装する。他の例では、コンピューティングシステム 110 および TTS システム 120 は、別個であり、互いに物理的に分離されている。コンピューティングシステム 120 は、分散システム (例えば、クラウドコンピューティング環境) を含み得る。

10

#### 【0028】

いくつかの実装形態では、ディープニューラルネットワーク 200 は、基準オーディオ信号 202 の大きなセットでトレーニングされる。各基準オーディオ信号 202 は、マイクロフォンによって録音された、韻律表現を有する人間音声の音声発話を含み得る。トレーニング中、ディープニューラルネットワーク 200 は、同じ音声発話に対して、異なる韻律を有する複数の基準オーディオ信号 202 を受信することができる (すなわち、同じ発話を複数の異なる方法で話すことができる)。ここで、基準オーディオ信号 202 は、コンテンツが同じであっても、音声発話の持続時間が変化するように可変長である。ディープニューラルネットワーク 200 は、各基準オーディオ信号 202 に関連付けられた韻律表現を、対応する固定長の発話埋め込み 204 に符号化/圧縮するように構成されている。ディープニューラルネットワーク 200 は、各固定長発話埋め込み 204 を、発話埋め込み 204 に関連付けられた基準オーディオ信号 202 の対応するトランスクリプト 206 とともに、(例えば、コンピューティングシステム 110 のメモリハードウェア 114 上の) 発話埋め込みストレージ 130 に格納することができる。ディープニューラルネットワーク 200 は、トランスクリプト 206 から抽出された言語的特徴を条件とする固定長発話埋め込み 204 を逆伝播して、各音節のピッチ、エネルギー、および持続時間の固定長フレームを生成することによってさらにトレーニングされ得る。

20

30

#### 【0029】

推論中、コンピューティングシステム 110 は、韻律モデル 300 を使用して、テキスト発話 310 の韻律表現 302 を予測することができる。韻律モデル 300 は、テキスト発話 310 のための発話埋め込み 204 を選択することができる。発話埋め込み 204 は、テキスト発話 310 の意図された韻律を表す。以下、図 2A ~ 図 2C および図 3A ~ 図 3D を参照してより詳細に説明されるように、韻律モデル 300 は、選択された発話埋め込み 204 を使用して、テキスト発話 310 の韻律表現 302 を予測することができる。韻律表現 302 は、テキスト発話 310 の予測されたピッチ、予測されたタイミング、および予測されたラウドネス (例えば、エネルギー) を含み得る。示されている例では、TTS システム 120 は、韻律表現 302 を使用して、テキスト発話 310 から意図された韻律を有する合成された音声 122 を生成する。

40

#### 【0030】

図 2A ~ 図 2C は、韻律の制御可能なモデルを提供するクロックワーク階層型変分オートエンコーダ (CHiVE) 300 (「オートエンコーダ 300」) の階層的言語構造 (例えば、図 1 のディープニューラルネットワーク) 200 を示す。韻律の制御可能なモデルは、所与の入力テキストの各音節について、所与の入力テキストからの固有のマッピングまたは他の言語仕様に依存することなく、音節の持続時間、音節のピッチ (F0) およびエネルギー (C0) 輪郭を共同で予測して、意図された/選択された韻律を有する合成音声 122 を生成することができる。オートエンコーダ 300 は、基準オーディオ信号 202 からサンプリングされた複数の固定長基準フレーム 210 を固定長発話埋め込み 20

50

4に符号化するエンコーダ部分320(図2A)と、固定長発話埋め込み204を復号化する方法を学習するデコーダ部分330(図2Bおよび図2C)とを含む。デコーダ部分330は、固定長発話埋め込み204を複数の固定長予測フレーム(fixed-length predicted frame)260に復号化することができる(例えば、発話埋め込み204に対してピッチ(F0)、エネルギー(C0)、またはスペクトル特性を予測するため)。これから明らかになるように、オートエンコーダ300は、デコーダ部分330から出力される予測フレーム260の数が、エンコーダ部分320に入力される基準フレーム210の数と等しくなるようにトレーニングされる。さらに、オートエンコーダ300は、基準フレーム210および予測フレーム260に関連付けられたデータが互いに実質的に一致するようにトレーニングされる。

10

#### 【0031】

図2Aを参照すると、エンコーダ部分320は、入力基準オーディオ信号202から固定長基準フレーム210のシーケンスを受信する。入力基準オーディオ信号202は、ターゲット韻律を含む、マイクロフォンによって記録された人間音声の音声発話を含み得る。エンコーダ部分320は、同じ音声発話に対して、異なる韻律を有する複数の基準オーディオ信号202を受信することができる(すなわち、同じ発話を複数の異なる方法で話すことができる)。例えば、同じ音声発話でも、話された基準が質問への回答である場合は、音声発話が質問である場合と比較して、韻律が異なり得る。基準フレーム210は、各々、5ミリ秒(ms)の持続時間を含み、基準オーディオ信号202に対するピッチの輪郭(F0)またはエネルギーの輪郭(C0)のうちの1つを表すことができる。並行して、エンコーダ部分320は、各々が5ミリ秒の持続時間を含み、基準オーディオ信号202に対するピッチの輪郭(F0)またはエネルギーの輪郭(C0)のうちの他の1つを表す基準フレーム210の第2のシーケンスを受信してもよい。したがって、基準オーディオ信号202からサンプリングされたシーケンス基準フレーム210は、持続時間、ピッチ輪郭、および/またはエネルギー輪郭を提供して、基準オーディオ信号202に対する韻律を表す。基準オーディオ信号202の長さまたは持続時間は、基準フレーム210の総数の合計と相関している。

20

#### 【0032】

エンコーダ部分320は、互いに対して相対的にクロックする、基準オーディオ信号202に対する基準フレーム210、音素220、220a、音節230、230a、単語240、240a、および文250、250aの階層レベルを含む。例えば、基準フレーム210のシーケンスに関連付けられたレベルは、音素220のシーケンスに関連付けられた次のレベルよりも速くクロックする。同様に、音節のシーケンス230に関連付けられたレベルは、音素330のシーケンスに関連付けられたレベルよりも遅く、かつ単語のシーケンス240に関連付けられたレベルよりも速くクロックする。したがって、より遅くクロックする層は、入力として、より速くクロックする層からの出力を受け取り、その結果、より速い層の最終クロック(すなわち、状態)の後の出力は、対応するより遅い層への入力として取られ、本質的にシーケンス間(sequence-to-sequence)エンコーダを提供する。示されている例では、階層レベルは、長短期記憶(LSTM)レベルを含む。

30

#### 【0033】

図2Aは、基準オーディオ信号202に対する階層レベルの例を示している。この例では、基準オーディオ信号202は、3つの単語240、240A~Cを有する1つの文240、240Aを含む。第1の単語240、240Aは、2つの音節230、230Aa~Abを含む。第2の単語240、240Bは、1つの音節230、230Baを含む。第3の単語240、240aは、2つの音節230、230Ca~Cbを含む。第1の単語240、240Aの第1の音節230、230Aaは、2つの音素220、220Aa1~Aa2を含む。第1の単語240、240Aの第2の音節230、230Abは、1つの音素220、220Ab1を含む。第2の単語240、240Bの第1の音節230、230Baは、3つの音素220、220Ba1~Ba3を含む。第3の単語240、

40

50

2 4 0 C の第 1 の音節 2 3 0 , 2 3 0 C a は、1 つの音素 2 2 0 , 2 2 0 C a 1 を含む。第 3 の単語 2 4 0 , 2 4 0 C の第 2 の音節 2 3 0 , 2 3 0 C b は、2 つの音素 2 2 0 , 2 2 0 C b 1 ~ C b 2 を含む。

#### 【 0 0 3 4 】

いくつかの例では、エンコーダ部分 3 2 0 は、最初に、基準フレーム 2 1 0 のシーケンスをフレームベースの音節埋め込み 2 3 2 , 2 3 2 A a ~ C b に符号化する。いくつかの実装形態では、基準フレーム 2 1 0 は、音素 2 2 0 A a 1 ~ 2 2 0 C b 2 のシーケンスを定義する。ここで、基準フレーム 2 1 0 のサブセットを 1 つまたは複数の音素 2 2 0 に符号化する代わりに、エンコーダ部分 3 2 0 は、代わりに、単音 ( p h o n e ) レベルの言語的特徴 2 2 2 , 2 2 2 A a 1 ~ C b 2 を単音特徴ベースの音節埋め込み 2 3 4 , 2 3 4 A a ~ C b に符号化することによって音素 2 2 0 を説明する。それぞれの音節埋め込み 2 3 2 , 2 3 4 は、対応する音節 2 3 0 に関連付けられた持続時間、ピッチ ( F 0 )、および/またはエネルギー ( C 0 ) を示す数値ベクトルを参照することができる。さらに、各音節埋め込み 2 3 2 , 2 3 4 は、音節 2 3 0 のレベルに対して対応する状態を示している。

#### 【 0 0 3 5 】

図 2 A を参照すると、階層的な層の斜めのハッチングパターンを含むブロックは、階層の特定のレベルの言語的特徴に対応する。フレームベースの音節埋め込み 2 3 2 および単音特徴ベースの音節埋め込み 2 3 4 を用いて、エンコーダ部分 3 2 0 は、これらの音節埋め込み 2 3 2 , 2 3 4 を他の言語的特徴とともに符号化する。例えば、エンコーダ部分 3 2 0 は、音節レベルの言語的特徴 2 3 6 , 2 3 6 A a ~ C b、単語レベルの言語的特徴 2 4 2 , 2 4 2 A ~ C、および/または文レベルの言語的特徴 2 5 2 , 2 5 2 A を用いて、音節埋め込み 2 3 2 , 2 3 4 を符号化する。言語的特徴 2 3 6 , 2 4 2 , 2 5 2 を用いて音節埋め込み 2 3 2 , 2 3 4 を符号化することにより、エンコーダ部分 3 2 0 は、基準オーディオ信号 2 0 2 のための発話埋め込み 2 0 4 を生成する。発話埋め込み 2 0 4 は、基準オーディオ信号 2 0 4 のそれぞれのトランスクリプト 2 0 6 (例えば、テキスト表現) と共にデータストレージ 1 3 0 (図 1) に格納され得る。トランスクリプト 2 0 6 から、言語的特徴 2 2 2 , 2 3 6 , 2 4 2 , 2 5 2 を抽出して、階層的言語構造 2 0 0 のトレーニングを調整する際に使用するために格納することができる。言語的特徴 (例えば、言語的特徴 2 2 2 , 2 3 6 , 2 4 2 , 2 5 2) は、各音素の個々の音、各音節が強勢を有するかまたは無強勢であるか、発話における各単語のタイプ (例えば、名詞/形容詞/動詞) および/または単語の位置、ならびに発話が質問であるかまたはフレーズであるかを含み得るが、これらに限定されない。

#### 【 0 0 3 6 】

図 2 A の例では、符号化ブロック 3 2 2 , 3 2 2 A a ~ C b は、言語的特徴 2 3 6 , 2 4 2 , 2 5 2 と音節埋め込み 2 3 2 , 2 3 4 との組み合わせを描写するために示されている。ここで、ブロック 3 2 2 は、音節のレートでシーケンス符号化されて、発話埋め込み 2 0 4 を生成する。例として、第 1 のブロック 3 2 2 A a は、第 2 のブロック 3 2 2 A b への入力として供給される。第 2 のブロック 3 2 2 A b は、第 3 のブロック 3 2 2 B a への入力として供給される。第 3 のブロック 3 2 2 C a は、第 4 のブロック 3 2 2 C a への入力として供給される。第 4 のブロック 3 2 2 C a は、第 5 のブロック 3 2 2 C b に供給される。いくつかの構成では、発話埋め込み 2 0 4 は、各基準オーディオ信号 2 0 2 の平均  $\mu$  および標準偏差 を含み、ここで、平均  $\mu$  および標準偏差 は、複数の基準オーディオ信号 2 0 2 のトレーニングデータに関するものである。

#### 【 0 0 3 7 】

いくつかの実装形態では、各音節 2 3 0 は、入力として、基準フレーム 2 1 0 のサブセットの対応するエンコーディングを受け取り、符号化されたサブセットにおける基準フレーム 2 1 0 の数に等しい持続時間を含む。示されている例では、最初の 7 つの固定長基準フレーム 2 1 0 は、音節 2 3 0 A a に符号化され、次の 4 つの固定長基準フレーム 2 1 0 は、音節 2 3 0 A b に符号化され、次の 1 1 個の固定長基準フレーム 2 1 0 は、音節 2 3 0 B a に符号化され、次の 3 つの固定長基準フレーム 2 1 0 は、音節 2 3 0 C a に符号化

10

20

30

40

50

され、最後の6つの固定長基準フレーム210は、音節230Cbに符号化される。したがって、音節230のシーケンスの各音節230は、音節230に符号化された基準フレーム210の数および対応するピッチおよび/またはエネルギー輪郭に基づく対応する持続時間を含むことができる。例えば、音節230Aaは、35ミリ秒に等しい持続時間を含み(すなわち、各々が5ミリ秒の固定長を有する6つの基準フレーム210)、音節230Abは、20ミリ秒に等しい持続時間を含む(すなわち、各々が5ミリ秒の固定長を有する4つの基準フレーム210)。したがって、基準フレーム210のレベルは、音節230のレベルにおける音節230Aaと次の音節230Abとの間の単一のクロッキングに対して合計10回クロックする。音節230の持続時間は、音節230のタイミング、および隣接する音節230の間の休止を示すことができる。

10

#### 【0038】

いくつかの実装形態では、エンコーダ部分320によって生成された発話埋め込み204は、基準オーディオ信号202の韻律を表す数値ベクトルを含む固定長の発話埋め込み204である。いくつかの例では、固定長発話埋め込み204は、「128」または「256」に等しい値を有する数値ベクトルを含む。エンコーダ部分320は、各々が同じ音声発話/フレーズに対応するが、異なる韻律を有する複数の基準オーディオ信号202を符号化することができ、すなわち、各基準オーディオ信号202は、同じ発話を伝達するが、異なって話される。

#### 【0039】

図2Bおよび図2Cを参照すると、変分オートエンコーダ300のデコーダ部分330は、発話の韻律を表す発話埋め込み204を最初に復号化することによって、複数の音節埋め込み232, 234(例えば、固定長音節埋め込み)を生成するように構成されている。トレーニング中、発話埋め込み204は、基準オーディオ信号202からサンプリングされた複数の固定長基準フレーム210を符号化することによって、図2Aのエンコーダ部分320から出力された発話埋め込み204を含み得る。したがって、デコーダ部分330は、トレーニング中に発話埋め込み204を逆伝播して、複数の固定長基準フレーム210に厳密に一致する複数の固定長予測フレーム260を生成するように構成されている。例えば、ピッチ(F0)およびエネルギー(C0)の両方に対する固定長予測フレーム260は、トレーニングデータとしてエンコーダ部分320に入力される基準オーディオ信号202の基準韻律と実質的に一致するターゲット韻律(例えば、予測された韻律)を表すために並行して生成されてもよい。追加的または代替的に、固定長予測フレーム260は、TTSシステム120(図1)に提供され得るスペクトルフレーム(例えば、メルフレーム)であってよい。いくつかの例では、TTSシステム120(図1)は、固定長予測フレーム260を使用して、固定長発話埋め込み204に基づいて、選択された韻律を有する合成音声122を生成する。例えば、TTSシステム120のユニット選択モジュール、WaveNetモジュール、またはニューラルボコーダは、フレーム260を使用して、意図された韻律を有する合成音声132を生成することができる。

20

30

#### 【0040】

示されている例では、デコーダ部分330は、エンコーダ部分320(図2A)から受信された発話埋め込み204(例えば、「256」または「128」の数値)を階層レベルに復号化する。例えば、階層レベルは、文250, 250b、単語240, 240b、音節230, 230b、音素220, 220b、および固定長予測フレーム260に対応するレベルを含む。具体的には、固定長発話埋め込み204は、デコーダ部分330に対する階層入力データの変分層に対応し、スタックされた階層レベルの各々は、階層入力データの長さに変動的にクロックされる長短期記憶(LSTM)処理ブロックを含む。例えば、音節レベル230は、単語レベル240よりも速く、かつ音素レベル220よりも遅くクロックする。各レベルの長方形のブロックは、それぞれの文、単語、音節、音素、またはフレームの1つまたは複数のLSTM処理セルに対応する。有利には、オートエンコーダ300は、単語レベル240のLSTM処理セルに最後の100単語にわたるメモリを与え、音節レベル230のLSTMセルに最後の100音節にわたるメモリを与え、音

40

50

素レベル 2 2 0 の L S T M セルに最後の 1 0 0 音素にわたるメモリを与え、固定長ピッチおよび/またはエネルギーフレーム 2 6 0 の L S T M セルに最後の 1 0 0 個の固定長フレーム 2 6 0 にわたるメモリを与える。固定長フレーム 2 6 0 が、それぞれ 5 ミリ秒の持続時間（例えば、フレームレート）を含む場合、対応する L S T M 処理セルは、最後の 5 0 0 ミリ秒（例えば、0.5 秒）にわたってメモリを提供する。

#### 【 0 0 4 1 】

図 2 B および図 2 C を参照すると、いくつかの例では、階層的言語構造 2 0 0 のデコーダ部分 3 3 0 は、エンコーダ部分 3 2 0 によって符号化された固定長発話埋め込み 2 0 4 を逆伝播する。例えば、図 2 B は、階層的言語構造 2 0 0 のデコーダ部分 3 3 0 が、固定長発話埋め込み 2 0 4 を、1 つの文 2 5 0 , 2 5 0 A , 3 つの単語 2 4 0 A ~ 2 4 0 C のシーケンス、5 つの音節 2 3 0 A a ~ 2 3 0 C b のシーケンス、および 9 つの音素 2 3 0 A a 1 ~ 2 3 0 C b 2 のシーケンスに逆伝播して、予測固定長フレーム 2 6 0 のシーケンスを生成することを示している。図 2 C のようないくつかの実装形態では、デコーダ部分 3 3 0 は、音節レベル 2 3 0 から音素レベル 2 2 0 に逆伝播するのではなく、音節レベル 2 3 0 からフレームレベル 2 1 0 に逆伝播する。ここで、このアプローチは、音素アライメントがモデル 2 0 0 にとって問題となる可能性がある状況において、予測されるフレーム 2 6 0 の精度を高めることができる。デコーダ部分 3 3 0 は、入力テキストの言語的特徴（例えば、言語的特徴 2 2 2 , 2 3 6 , 2 4 2 , 2 5 2 ）に基づいて条件付けられる。より速くクロックする層からの出力が、より遅くクロックする層によって入力として受け取られ得るエンコーダ部分 3 2 0（例えば、図 2 A に示されるような）とは対照的に、デコーダ部分 3 3 0 は、より速くクロックする層に供給する、より遅くクロックする層からの出力を含み、これにより、より遅くクロックする層の出力が、それに付加されたタイミング信号を用いて、各クロックサイクルで、より速くクロックする層の入力に分配される。

#### 【 0 0 4 2 】

図 3 A ~ 図 3 C を参照すると、いくつかの実装形態では、オートエンコーダ 3 0 0 は、階層的言語構造 2 0 0 を使用して、推論中に所与のテキスト発話 3 1 0 の韻律表現を予測する。例えば、オートエンコーダ 3 0 0 は、所与のテキスト発話 3 1 0 の各音節 2 3 0 に対して、音節 2 3 0 の持続時間およびピッチ F 0 および/またはエネルギー C 0 輪郭を一緒に予測することによって、所与のテキスト発話 3 1 0 に対する韻律表現を予測する。テキスト発話 3 1 0 は、テキスト発話 3 1 0 の適切な韻律を示すためのコンテキスト、セマンティック情報、または語用論情報を提供しないので、オートエンコーダ 3 0 0 は、潜在変数として発話埋め込み 2 0 4 を選択して、テキスト発話 3 1 0 の意図された韻律を表す。

#### 【 0 0 4 3 】

発話埋め込み 2 0 4 は、発話埋め込みデータストレージ 1 3 0（図 1）から選択することができる。ストレージ 1 3 0 内の各発話埋め込み 2 0 4 は、トレーニング中に、対応する可変長基準オーディオ信号 2 0 2（図 2 A）からエンコーダ部分 3 2 0（図 2 A）によって符号化され得る。具体的には、エンコーダ部分 3 1 0 は、トレーニング中に可変長基準オーディオ信号 2 0 2 の韻律を、固定長発話埋め込み 2 0 4 に圧縮し、推論時にデコーダ部分 3 3 0 によって使用するため、各発話埋め込み 2 0 4 を、対応する基準オーディオ信号 2 0 2 のトランスクリプト 2 0 6 と共に、発話埋め込みデータストレージ 1 3 0 に格納する。示される例では、オートエンコーダ 3 0 0 は、最初に、テキスト発話 3 1 0 に厳密に一致するトランスクリプト 2 0 6 を有する発話埋め込み 2 0 4 を見つけ、次に、発話埋め込み 2 0 4 のうちの 1 つを選択して、所与のテキスト発話 3 1 0 の韻律表現 3 0 2（図 1）を予測することができる。いくつかの例では、固定長の発話埋め込み 2 0 4 は、ターゲット韻律の特定のセマンティクスおよび語用論を表す可能性が高い埋め込み 2 0 4 の潜在空間内の特定の点を選ぶことによって選択される。他の例では、潜在空間は、テキスト発話 3 1 0 の意図された韻律を表すために、ランダムな発話埋め込み 2 0 4 を選択するためにサンプリングされる。さらに別の例では、オートエンコーダ 3 0 0 は、テキスト発話 3 1 0 の言語的特徴に対する最も可能性の高い韻律を表すために、厳密に一致するトラ

10

20

30

40

50

ンスクリプト 206 を有する発話埋め込み 204 の平均を選択することによって、潜在空間を多次元ユニットガウシアンとしてモデル化する。例えば、オートエンコーダ 300 は、発話埋め込み 204 を選択してテキスト発話 310 に対する韻律表現 302 を生成する場合に、各発話埋め込み 204 に関連付けられた平均  $\mu$  および  $\sigma$  / または標準偏差  $\sigma$  を使用する。トレーニングデータの韻律の変化が適度中立である場合、発話埋め込み 204 の平均を選択する最後の例は妥当な選択である。

#### 【0044】

図 3A ~ 図 3C は、階層的言語構造 200 の単語レベル 240 で表される 3 つの単語 240A, 240B, 240C を有するテキスト発話 310 を示している。第 1 の単語 240A は、音節 230Aa, 230Ab を含む。第 2 の単語 240B は、1 音節 230Ba を含む。第 3 の単語 240C は、音節 230Ca, 230Cb を含む。したがって、階層的言語構造 200 の音節レベル 230 は、テキスト発話 310 の 5 つの音節 230Aa ~ 230Cb のシーケンスを含む。LSTM 処理セルの音節レベル 230 において、オートエンコーダ 300 は、第 1 の音節埋め込み (例えば、フレームベースの音節埋め込み 232) および第 2 の音節埋め込み (例えば、単音特徴ベースの音節埋め込み 234) を生成 / 出力するように構成されている。第 1 の音節埋め込み 232Aa, 232Ab, 232Ba, 232Ca, および 232Cb を出力するために、オートエンコーダ 300 は、以下の入力、すなわち、固定長発話埋め込み 204、テキスト発話 310 に関連付けられた発話レベルの言語的特徴 (例えば、文レベルの言語的特徴 252)、音節 230 を含む単語 240 に関連付けられた単語レベルの言語的特徴 242、音節 230 に対する音節レベルの言語的特徴 236 を使用する。第 2 の音節埋め込み 234Aa, 234Ab, 234Ba, 234Ca, および 234Cb を出力するために、オートエンコーダ 300 は、音節 230 に関連付けられた単音レベルの言語的特徴 222 を使用する。発話レベルの言語的特徴 252 は、テキスト発話 320 が質問であるかどうか、質問への回答であるかどうか、フレーズであるかどうか、文であるかどうかなどを含むがこれらに限定されない。いくつかの例では、DVector または他の複雑な話者アイデンティティ表現が、発話レベルの言語的特徴 252 として含まれ得る。単語レベルの言語的特徴 242 は、これらに限定されないが、単語タイプ (例えば、名詞、代名詞、動詞、形容詞、副詞など) およびテキスト発話 310 における単語の位置を含み得る。音節レベルの言語的特徴 236 は、これに限定されないが、音節 230 が強勢を有するか無強勢であるかを含み得る。いくつかの実装形態では、DVector または他の複雑な話者アイデンティティ表現が、発話レベル (例えば、単語レベル 240 または音節レベル 230) より下のレベルにおける言語的特徴として含まれ得る。

#### 【0045】

示されている例では、音節レベル 230 内の各音節 230Aa, 230Ab, 230Ba, 230Ca, 230Cb は、個々の固定長予測ピッチ ( $F_0$ ) フレーム 260, 260F0 (図 3A) を復号化するため、個々の固定長予測エネルギー ( $C_0$ ) フレーム 260, 260C0 (図 3B) を復号化するため、および / または個々の固定長スペクトル ( $M_0$ ) フレーム 260, 260M0 を復号化するため、対応する音節埋め込み 232Aa ~ Cb, 234Aa ~ Cb を出力する対応する LSTM 処理セルに関連付けられ得る。いくつかの実装形態では、オートエンコーダ 300 は、2 つ以上のタイプのフレーム  $F_0$ ,  $C_0$ ,  $M_0$  を並列に復号化する。図 3A は、音節 230 に対する持続時間 (タイミングおよび休止) およびピッチ輪郭を示す複数の固定長予測ピッチ ( $F_0$ ) フレーム 260F0 を含む、音節レベル 230 における各音節 230 を示している。ここで、持続時間およびピッチ輪郭は、音節 230 の韻律表現に対応する。図 3B は、音節 230 に対する持続時間およびエネルギー輪郭を示す複数の固定長予測エネルギー ( $C_0$ ) フレーム 260C0 を含む、音節レベル 240 における各音節 230 を示している。

#### 【0046】

音節レベル 230 の第 1 の音節 230Aa (すなわち、LSTM 処理セル Aa) は、対応する音節埋め込み 232Aa を生成するための入力として、固定長発話埋め込み 204

10

20

30

40

50

、テキスト発話 3 1 0 に関連付けられた発話レベルの言語的特徴 2 5 2、第 1 の単語 2 3 0 A に関連付けられた単語レベルの言語的特徴 2 4 2 A、および音節 2 3 0 A a に対する音節レベルの言語的特徴 2 3 6 A a を受け取る。音節レベル 2 3 0 の第 2 の音節 2 3 0 A b は、対応する音節埋め込み 2 3 2 A b を生成するための入力として、固定長発話埋め込み 2 0 4、テキスト発話 3 1 0 に関連付けられた発話レベルの言語的特徴 2 5 2、第 1 の単語 2 4 0 A に関連付けられた単語レベルの言語的特徴 2 4 2 A、および音節 2 3 0 A b に対する対応する音節レベルの言語的特徴 2 3 6 を受け取る。ここで、オートエンコーダ 3 0 0 がテキスト発話 3 1 0 を音節埋め込み 2 3 2、2 3 4 に復号化する場合、これらの埋め込み 2 3 2、2 3 4 は、各音節 2 3 0 に対して状態 3 3 2 (例えば、状態 3 3 2、3 3 2 A a ~ C b として示される) を形成する。オートエンコーダ 3 0 0 が後続の状態 3 3 2 を形成する場合、オートエンコーダ 3 0 0 は、先行する音節 2 3 0 の状態も受け取る。換言すると、音節レベル 2 3 0 の各 L S T M 処理セルは、音節レベル 2 3 0 の直前の L S T M 処理セルの状態 2 3 2 を受け取る。いくつかの構成では、音節レベル 2 3 0 の各 L S T M 処理セルは、現在の音節レベル 2 3 0 の L S T M 処理セルに先行する各状態 2 3 2 を受け取る。例えば、第 2 の音節 2 3 0 A b に関連付けられた L S T M 処理セルは、先行する第 1 の音節 2 3 0 A a の状態 3 3 2、3 3 2 A a を受け取る。図 3 A ~ 図 3 C に示されるように、オートエンコーダ 3 0 0 は、同様の方法で、音節レベル 2 3 0 における音節 2 3 0 B a、2 3 0 C a、2 3 0 C b の残りのシーケンスに対して、対応する音節埋め込み 2 3 2 B a ~ C b、2 3 4 B a ~ C B を生成する。追加的または代替的に、オートエンコーダ 3 0 0 は、連結によって、1 つまたは複数の音節埋め込み 2 3 2、2 3 4 を用いて、より高いレベルの言語的特徴 (例えば、文の特徴 2 5 2、単語の特徴 2 4 2 など) を復号化することができる。

10

20

#### 【 0 0 4 7 】

図 3 D を参照すると、階層的言語構造 2 0 0 の音素レベル 2 2 0 は、9 つの音素 2 2 0 A a 1 ~ 2 2 0 C b 2 のシーケンスを含む。いくつかの実装形態では、オートエンコーダ 3 0 0 は、各音素 2 2 0 A a 1 ~ 2 2 0 C b 2 に関連付けられた音素レベルの言語的特徴 2 2 2 を、音素符号化状態 2 2 4、2 2 4 A a 1 ~ C a 2 に符号化する。音素レベルの言語的特徴 2 2 2 は、対応する音素 2 2 0 の音のアイデンティティを含み得るがこれに限定されない。図 3 D に示されるように、アテンション機構 3 4 0 は、符号化された音素状態 2 2 4 を形成する単音レベルの特徴 2 2 2、2 2 A a 1 ~ C b 2 にアテンションを与えるように構成されている。このアプローチでは、オートエンコーダ 3 0 0 は、音素持続時間に依存する必要もなく、音素持続時間を予測もする必要もなく、むしろ、単音レベルの言語的特徴 2 2 2 へのアテンションを有するアテンション機構 3 4 0 を使用して、各音節 2 3 0 の音節持続時間 2 3 8 を予測する。

30

#### 【 0 0 4 8 】

一般的に言えば、アテンション機構 3 4 0 は、入力を出力と (例えば、スコアリングすることによって) 相関させるアライメントモデルである。符号化された隠された状態では、アテンション機構 3 4 0 は、各出力 (例えば、予測されたフレーム 2 6 0) について、各隠された状態 (例えば、符号化された音素状態 2 2 4) がどれだけのアテンションが考慮されるべきかを定義する重みのセットを形成し得る。異なるアラインメントスコア関数を使用する異なるタイプのアテンション機構 3 4 0 があってもよい。これらのアテンション機構 3 4 0 のいくつかの例は、コンテンツベースのアテンション、加法 (a d d i t i v e) アテンション、場所ベースのアテンション、一般的 (g e n e r a l) アテンション、ドット積 (d o t - p r o d u c t) アテンション、およびスケールされたドット積アテンションを含む。アテンション機構 3 4 0 のより広いカテゴリーは、セルフアテンション、グローバル/ソフトアテンション、および/またはローカル/ハードアテンションを含む。いくつかの例では、アテンション機構 3 4 0 のアライメントスコアは、単一の隠れ層 (例えば、符号化された音素状態 2 2 4 の音素層 2 2 0) を備えたフィードフォワードネットワークによってパラメータ化される。これらの例では、フィードフォワードネットワークは、オートエンコーダ 3 0 0 の他の部分と共同でトレーニングされ得る。

40

50

## 【 0 0 4 9 】

オートエンコーダ 3 0 0 のアテンション機構 3 4 0 は、これらのタイプのアテンションモデルの任意のものに基づくことができる。換言すると、オートエンコーダ 3 0 0 は、異なるアテンションアプローチに従って、単音レベルの言語的特徴 2 2 2 に基づいて隠された状態 2 2 4 をスコアリングするように設計され得る。いくつかの例では、アテンション機構 3 4 0 は、単調にシフトする、位置に敏感なアテンション（すなわち、位置ベースのアテンション）の形をとる。ここで、アテンション機構 3 0 0 の位置に敏感なアテンションは、ウィンドウによって制約され得る。いくつかの構成では、ウィンドウは、オートエンコーダ 3 0 0 によって処理されている現在の音節 2 3 0 の音素情報に対応する。他の構成では、ウィンドウは、オートエンコーダ 3 0 0 によって処理されている現在の音節 2 3 0 に隣接する音節 2 3 0（例えば、前または後）を含むように拡張され得る。いくつかの実装形態では、アテンション機構 3 4 0 は、セルフアテンション機構に基づくトランスフォーマー（transformer）である。アテンション機構 3 4 0 としてトランスフォーマーを使用する場合、アテンション機構 3 4 0 は、符号化された音素状態 2 2 4 における符号化された単音レベルの言語的特徴 2 2 2 ではなく、単音レベルの言語的特徴 2 2 2 を入力として使用し得る。

10

## 【 0 0 5 0 】

図 3 A ~ 図 3 C を引き続き参照すると、いくつかの例では、オートエンコーダ 3 0 0 は、2 つのデコーダを使用する（例えば、円で表されている）。オートエンコーダ 3 0 0 のこれらのデコーダの各々は、入力として音節埋め込み 2 3 2 , 2 3 4 を受け取ることができる。第 1 のデコーダ（例えば、第 1 の音節状態 3 3 2）は、アテンション機構 3 4 0 のアテンションとともに、音節埋め込み 2 3 2 , 2 3 4 を使用して、特定の音節 2 3 0 に対するフレーム数を予測する。次に、その特定の音節 2 3 0 について、オートエンコーダ 3 0 0 は、（例えば、第 2 の音節状態 3 3 4 の）第 2 のデコーダを使用して、（例えば、第 1 の音節状態 3 3 2 の）第 1 のデコーダによって予測されたフレーム数をシーケンス復号化する。予測フレーム 2 6 0 の数に基づいて、オートエンコーダ 3 0 0 は、予測フレーム 2 6 0 を生成する。換言すると、各音節 2 3 0 に対してアテンション機構 3 4 0 を使用して、予測音節持続時間は、固定長予測ピッチ（F 0）フレーム 2 6 0 F 0（例えば、図 3 A に示されるような）の数、固定長予測エネルギー（C 0）フレーム 2 6 0 C 0（例えば、図 3 B に示されるような）の数、および/または固定長予測スペクトル（M 0）フレーム 2 6 0 M 0（例えば、図 3 C に示されるような）の数を決定し得る。図 3 A ~ 図 3 C は単一のアテンション機構 3 4 0 を示しているが、これは、単に説明を簡略化するためである。テキスト発話 3 1 0 の各音節 2 3 0 について、オートエンコーダ 3 0 0 は、状態 3 3 2 および状態 3 3 4 でアテンション機構 3 4 0 を使用して、音節 2 3 0 に対応する予測フレーム 2 6 0 を生成する。例えば、図 3 A ~ 図 3 C は、点線のボックスと相互作用するアテンション機構 3 4 0 を示しており、オートエンコーダ 3 0 0 が、各音節 2 3 0 の各状態 3 3 2 , 3 3 4 でアテンション機構 3 4 0 を使用することを示している。

20

30

## 【 0 0 5 1 】

アテンション機構 3 4 0 および音節埋め込み 2 3 2 A a , 2 3 4 A a に基づいて、オートエンコーダ 3 0 0 は、7 つの予測フレーム 2 6 0 に対応する第 1 の音節 2 3 0 A a の音節持続時間 2 3 8 , 2 3 8 A a を予測する。第 1 の音節 2 3 0 A a に対する状態 3 3 4 , 3 3 4 A a で、オートエンコーダ 3 0 0 は、7 つの予測フレーム 2 6 0 を復号化する。第 2 の音節 2 3 0 A b について、アテンション機構 3 4 0 および音節埋め込み 2 3 2 A a ~ A b , 2 3 4 A a ~ A b に基づいて、オートエンコーダ 3 0 0 は、4 つの予測フレーム 2 6 0 の音節持続時間 2 3 8 , 2 3 8 A b を予測する。第 2 の音節 2 3 0 A b に対する状態 3 3 4 , 3 3 4 A b で、オートエンコーダ 3 0 0 は、4 つの予測フレーム 2 6 0 を復号化する。第 3 の音節 2 3 0 B a について、アテンション機構 3 4 0 および音節埋め込み 2 3 2 A a ~ B a , 2 3 4 A a ~ B a に基づいて、オートエンコーダ 3 0 0 は、1 1 個の予測フレーム 2 6 0 の音節持続時間 2 3 8 , 2 3 8 B a を予測する。第 3 の音節 2 3 0 B a に対する状態 3 3 4 , 3 3 4 B a で、オートエンコーダ 3 0 0 は、1 1 個の予測フレーム 2

40

50

60を復号化する。第4の音節230Caについて、アテンション機構340および音節埋め込み232Aa~Ca, 234Aa~Caに基づいて、オートエンコーダ300は、3つの予測フレーム260の音節持続時間238, 238Caを予測する。第4の音節230Caに対する状態334, 334Caで、オートエンコーダ300は、3つの予測フレーム260を復号化する。第5の音節230Cbについて、アテンション機構340および音節埋め込み232Aa~Cb, 234Aa~Cbに基づいて、オートエンコーダ300は、6つの予測フレーム260の音節持続時間238, 238Cbを予測する。第5の音節230Cbに対する状態334, 334Cbで、オートエンコーダ300は、6つの予測フレーム260を復号化する。

#### 【0052】

ここで、このシステム100は、(例えば、音素レベルが階層的言語構造200に直接影響を与えることなく)フレーム260を予測するために、オートエンコーダ300を使用して音節230に焦点を合わせる。しかしながら、本明細書のアプローチは、階層的言語構造200の他の層(例えば、文250、単語240など)に焦点を合わせるように適合されてもよい。例えば、オートエンコーダ300は、フレーム260を予測するために単語240または文250に焦点を合わせる。これらの他のアプローチでは、1つまたは複数の異なるレベル(例えば、音素レベルと同様の)を、特定のレベルに焦点を合わせるためのトレードオフとして、階層的言語構造200から削除することができる。これらのアプローチでは、アテンション機構340は、1つまたは複数の特定の言語層に焦点を合わせるように適宜適合されてもよい。

#### 【0053】

図4は、テキスト発話310の韻律表現302を予測する方法400のための動作の例示的な配置のフローチャートである。方法400は、図1~図3Cを参照して説明することができる。図1のコンピュータシステム110上にあるメモリハードウェア124は、データ処理ハードウェア112によって実行されると、データ処理ハードウェア112に方法400のための動作を実行させる命令を格納している。動作402において、方法400は、テキスト発話310を受信することを含む。テキスト発話310は、少なくとも1つの単語240を有し、各単語240は、少なくとも1つの音節230を有し、各音節230は、少なくとも1つの音素220を有する。動作404において、方法400は、テキスト発話310のための発話埋め込み204を選択することを含む。発話埋め込み204は、意図された韻律を表す。本明細書で使用される場合、選択された発話埋め込み204は、テキスト発話310から意図された韻律を有する合成音声122を生成するためにTTSシステム120が使用するための、テキスト発話310の韻律表現302を予測するために使用される。発話埋め込み204は、固定長の数値ベクトルによって表すことができる。数値ベクトルは、「256」に等しい値を含み得る。テキスト発話310に対して発話埋め込み204を選択するために、データ処理ハードウェア112は、最初にデータストレージ130に問い合わせて、テキスト発話310に厳密に一致するトランスクリプト206を有する発話埋め込み204を見つけ、次に、発話埋め込み204を選択して、所与のテキスト発話310の韻律表現302を予測することができる。いくつかの例では、固定長の発話埋め込み204は、ターゲット韻律の特定のセマンティクスおよび語用論を表す可能性が高い埋め込み204の潜在空間内の特定の点を選ぶことによって選択される。他の例では、潜在空間は、テキスト発話310の意図された韻律を表すために、ランダムな発話埋め込み204を選択するためにサンプリングされる。さらに別の例では、データ処理ハードウェア112は、テキスト発話310の言語的特徴に対する最も可能性の高い韻律を表すために、厳密に一致するトランスクリプト206を有する発話埋め込み204の平均を選択することによって、潜在空間を多次元ユニットガウシアンとしてモデル化する。

#### 【0054】

動作406において、各音節230について、選択された発話埋め込み204を使用して、方法400は、音節230の各音素220に対する言語的特徴222へのアテンシヨ

10

20

30

40

50

ン機構 340 によるアテンションに基づいて、音節 230 の韻律音節埋め込み 232, 234 を復号化することによって、音節 230 の持続時間 238 を予測することを含む。動作 408 において、各音節 230 について、選択された発話埋め込み 204 を使用して、方法 400 は、音節 230 に対して予測された持続時間 238 に基づいて、複数の固定長予測フレーム 260 を生成することを含む。

#### 【0055】

図 5 は、本文書に記載されているシステム（例えば、階層構造 200、オートエンコーダ 300、および/またはアテンション機構 340）および方法（例えば、方法 400）を実装するために使用することができる例示的なコンピューティングデバイス 500 の概略図である。コンピューティングデバイス 500 は、ラップトップ、デスクトップ、ワークステーション、パーソナルデジタルアシスタント、サーバ、ブレードサーバ、メインフレーム、および他の適切なコンピュータなどの様々な形態のデジタルコンピュータを表すことが意図されている。ここに示されているコンポーネント、それらの接続および関係、並びにそれらの機能は、単なる例示を意味するものであり、この文書で説明および/または主張されている発明の実施を制限することを意味するものではない。

#### 【0056】

コンピューティングデバイス 500 は、プロセッサ 510（例えば、データ処理ハードウェア）と、メモリ 520（例えば、メモリハードウェア）と、ストレージデバイス 530 と、メモリ 520 および高速拡張ポート 550 に接続する高速インタフェース/コントローラ 540 と、低速バス 570 およびストレージデバイス 530 に接続する低速インタフェース/コントローラ 560 とを含む。コンポーネント 510, 520, 530, 540, 550, および 560 の各々は、様々なバスを使用して相互接続されており、共通のマザーボードに、または必要に応じて他の方法で取り付けることができる。プロセッサ 510 は、高速インタフェース 540 に結合されたディスプレイ 580 などの外部入力/出力デバイスにグラフィカルユーザインタフェース（GUI）のためのグラフィカル情報を表示するために、メモリ 520 またはストレージデバイス 530 に格納された命令を含む、コンピューティングデバイス 500 内で実行するための命令を処理することができる。他の実装形態では、複数のメモリおよびメモリのタイプとともに、必要に応じて、複数のプロセッサおよび/または複数のバスを使用することができる。また、複数のコンピューティングデバイス 500 が接続されてもよく、各デバイスは、（例えば、サーババンク、ブレードサーバのグループ、またはマルチプロセッサシステムとして）必要な動作の一部を提供する。

#### 【0057】

メモリ 520 は、コンピューティングデバイス 500 内に非一時的に情報を格納する。メモリ 520 は、コンピュータ可読媒体、揮発性メモリユニット（複数可）、または不揮発性メモリユニット（複数可）であってよい。非一時的メモリ 520 は、コンピューティングデバイス 500 によって使用するために一時的または永続的にプログラム（例えば、命令のシーケンス）またはデータ（例えば、プログラム状態情報）を格納するために使用される物理デバイスであってよい。不揮発性メモリの例には、フラッシュメモリおよび読み取り専用メモリ（ROM）/プログラム可能読み取り専用メモリ（PROM）/消去可能プログラム可能読み取り専用メモリ（EPROM）/電子的に消去可能プログラム可能読み取り専用メモリ（EEPROM）（例えば、ブートプログラムなどのファームウェアに通常使用される）が含まれるが、これらに限定されない。揮発性メモリの例には、ランダムアクセスメモリ（RAM）、ダイナミックランダムアクセスメモリ（DRAM）、スタティックランダムアクセスメモリ（SRAM）、相変化メモリ（PCM）、およびディスクまたはテープが含まれるが、これらに限定されない。

#### 【0058】

ストレージデバイス 530 は、コンピューティングデバイス 500 に大容量ストレージデバイスを提供することができる。いくつかの実装形態では、ストレージデバイス 530 は、コンピュータ可読媒体である。様々な異なる実施形態では、ストレージデバイス 53

10

20

30

40

50

0 は、フロッピーディスク（登録商標）デバイス、ハードディスクデバイス、光ディスクデバイス、またはテープデバイス、フラッシュメモリまたは他の類似のソリッドステートストレージデバイス、またはストレージエリアネットワークまたは他の構成におけるデバイスを含むデバイスのアレイであってよい。追加の実装形態では、コンピュータプログラム製品は、情報担体に有体的に具現化される。コンピュータプログラム製品は、実行されると、上述したような1つまたは複数の方法を実行する命令を含む。情報担体は、メモリ520、ストレージデバイス530、またはプロセッサ510上のメモリなどの、コンピュータまたは機械で読み取り可能な媒体である。

#### 【0059】

高速コントローラ540は、コンピューティングデバイス500のための帯域幅集約的な動作を管理し、一方、低速コントローラ560は、より少ない帯域幅を消費する動作を管理する。このような職務の割り当ては例示に過ぎない。いくつかの実装形態では、高速コントローラ540は、メモリ520、ディスプレイ580（例えば、グラフィックプロセッサまたはアクセラレータを介して）に結合され、および様々な拡張カード（図示せず）を受け入れることができる高速拡張ポート550に結合される。いくつかの実装形態では、低速コントローラ560は、ストレージデバイス530および低速拡張ポート590に結合されている。様々な通信ポート（例えば、USB、ブルートゥース（登録商標）、イーサネット（登録商標）、ワイヤレスイーサネット）を含み得る低速拡張ポート590は、キーボード、ポインティングデバイス、スキャナなどの1つまたは複数の入力/出力デバイスに、またはスイッチやルーターなどのネットワークングデバイスに、例えば、ネットワークアダプタを介して結合され得る。

#### 【0060】

コンピューティングデバイス500は、図に示されるように、複数の異なる形態で実装することができる。例えば、それは、標準サーバ500aとして、またはそのようなサーバ500aのグループに複数回、ラップトップコンピュータ500bとして、またはラックサーバシステム500cの一部として実装することができる。

#### 【0061】

本明細書に記載のシステムおよび技術の様々な実装形態は、デジタル電子および/または光回路、集積回路、特別に設計されたASIC（特定用途向け集積回路）、コンピュータハードウェア、ファームウェア、ソフトウェア、および/またはそれらの組み合わせで実現することができる。これらの様々な実装形態は、ストレージシステム、少なくとも1つの入力デバイス、および少なくとも1つの出力デバイスからデータおよび命令を受信し、それらへデータおよび命令を送信するために結合された、専用または汎用であり得る少なくとも1つのプログラム可能なプロセッサを含むプログラム可能なシステム上で実行可能および/または解釈可能な1つまたは複数のコンピュータプログラムにおける実装形態を含むことができる。

#### 【0062】

これらのコンピュータプログラム（プログラム、ソフトウェア、ソフトウェアアプリケーション、またはコードとしても知られる）は、プログラム可能なプロセッサのマシン命令を含み、高レベルの手続き型および/またはオブジェクト指向のプログラミング言語、および/またはアセンブリ/マシン言語で実装することができる。本明細書で使用される場合、「機械可読媒体」および「コンピュータ可読媒体」という用語は、機械命令および/またはデータをプログラム可能なプロセッサに提供するために使用される任意のコンピュータプログラム製品、非一時的なコンピュータ可読媒体、装置および/またはデバイス（例えば、磁気ディスク、光ディスク、メモリ、プログラマブルロジックデバイス（PLD））を指し、機械命令を機械可読信号として受信する機械可読媒体を含む。「機械可読信号」という用語は、プログラム可能なプロセッサに機械命令および/またはデータを提供するために使用される任意の信号を指す。

#### 【0063】

本明細書に記載のプロセスおよび論理フローは、1つまたは複数のコンピュータプログ

10

20

30

40

50

ラムを実行して、入力データを操作し、出力を生成することによって機能を実行する1つまたは複数のプログラム可能なプロセッサによって実行することができる。プロセスおよびロジックフローは、例えば、FPGA（フィールドプログラマブルゲートアレイ）またはASIC（特定用途向け集積回路）などの特定用途のロジック回路によっても実行できる。コンピュータプログラムの実行に適したプロセッサは、例として、汎用マイクロプロセッサおよび専用マイクロプロセッサの両方、および任意の種類のデジタルコンピュータの任意の1つまたは複数のプロセッサを含む。一般に、プロセッサは、読み取り専用メモリまたはランダムアクセスメモリ、あるいはその両方から命令とデータを受け取る。コンピュータの重要な素子は、命令を実行するためのプロセッサと、命令およびデータを格納するための1つまたは複数のメモリデバイスである。一般に、コンピュータは、データを格納するための1つまたは複数の大容量ストレージデバイス、例えば、磁気、光磁気ディスク、または光ディスクを含むか、またはそれらからデータを受信し、または転送し、あるいはその両方のために動作可能に結合される。しかしながら、コンピュータがそのようなデバイスを有する必要はない。コンピュータプログラム命令およびデータを格納するのに適したコンピュータ可読媒体は、あらゆる形態の不揮発性メモリ、媒体およびメモリデバイスを含み、それらは、例えば、半導体メモリデバイス、例えば、EPROM、EEPROM、およびフラッシュメモリデバイス、磁気ディスク、例えば、内蔵ハードディスクまたはリムーバブルディスク、光磁気ディスク、およびCD-ROMおよびDVD-ROMディスクを含む。プロセッサおよびメモリは、特定用途の論理回路によって補完されてよく、または特定用途の論理回路に組み込まれてもよい。

10

20

#### 【0064】

ユーザとのインタラクションを提供するために、本開示の1つまたは複数の態様は、ユーザに情報を表示するためのディスプレイデバイス、例えば、CRT（陰極線管）、LCD（液晶ディスプレイ）モニタ、またはタッチスクリーンを有するコンピュータ上に実装されてもよく、コンピュータは、任意選択で、キーボードおよびポインティングデバイス、例えば、マウスまたはトラックボールを有し、これによって、ユーザがコンピュータに入力を提供できる。他の種類のデバイスを使用して、ユーザとのインタラクションを提供することもでき、例えば、ユーザに提供されるフィードバックは、例えば、視覚的フィードバック、聴覚的フィードバック、または触覚的フィードバックなど、任意の形態の感覚的フィードバックであってよく、ユーザからの入力は、音響、音声、または触覚入力を含む任意の形式で受け取ることができる。加えて、コンピュータは、ユーザによって使用されるデバイスとの間でドキュメントを送受信することにより、ユーザと対話（interact）でき、これは、例えば、Webブラウザから受信された要求に応答して、ユーザのクライアントデバイス上のWebブラウザにWebページを送信することによって行われる。

30

#### 【0065】

複数の実装形態が説明された。それにもかかわらず、本開示の技術思想および範囲から逸脱することなく、様々な修正を行うことができることが理解されるであろう。したがって、他の実装形態は、以下の特許請求の範囲内にある。

40

50

【図面】  
【図 1】

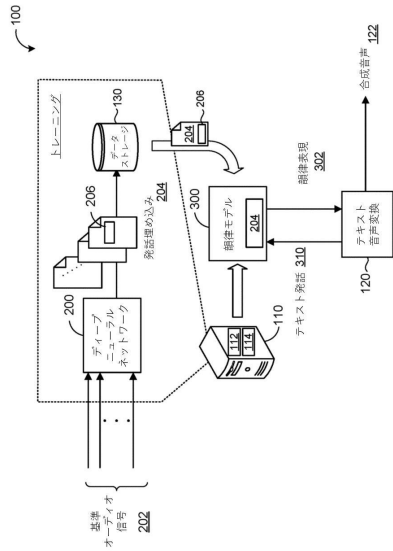


FIG. 1

【図 2 A】

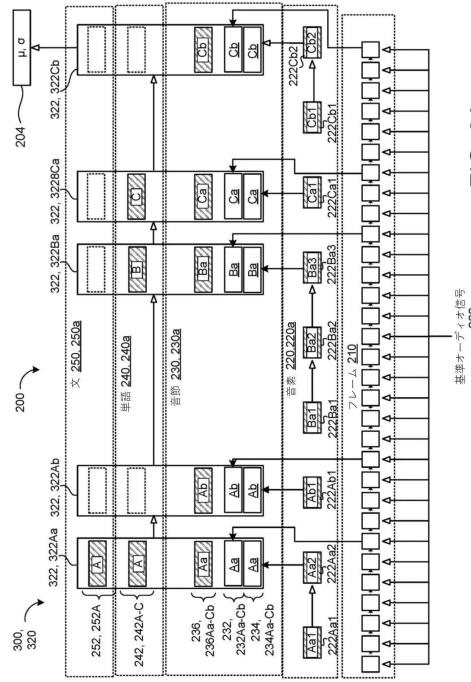


FIG. 2A

10

20

【図 2 B】

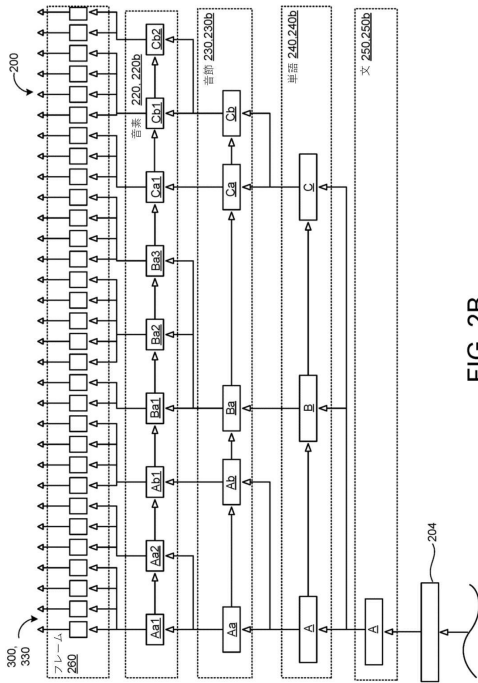


FIG. 2B

【図 2 C】

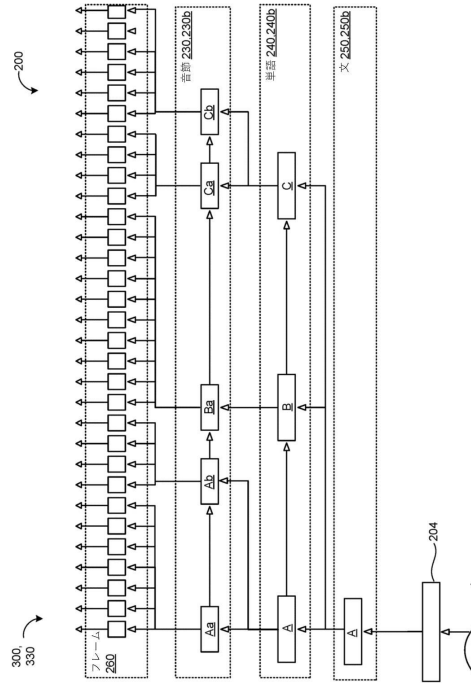


FIG. 2C

30

40

50

【図 3 A】

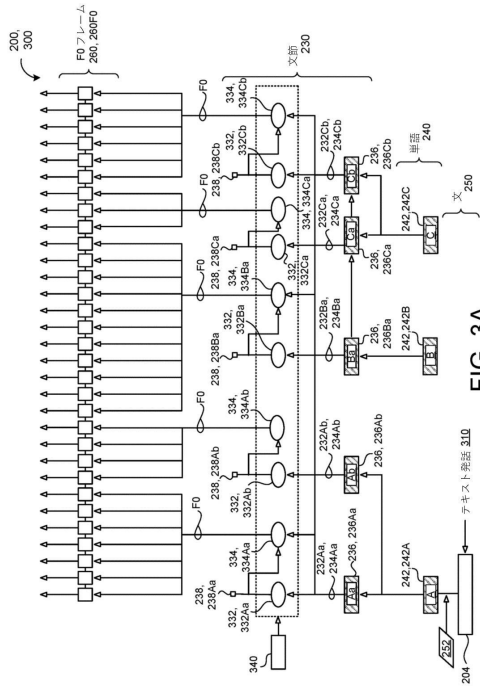


FIG. 3A

【図 3 B】

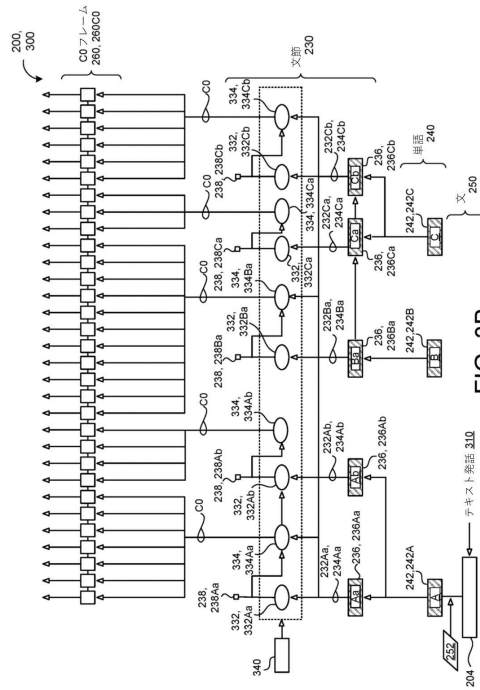


FIG. 3B

【図 3 C】

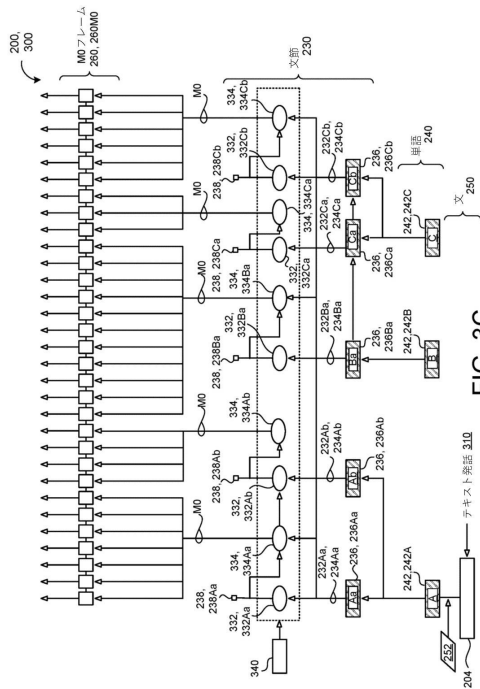


FIG. 3C

【図 3 D】

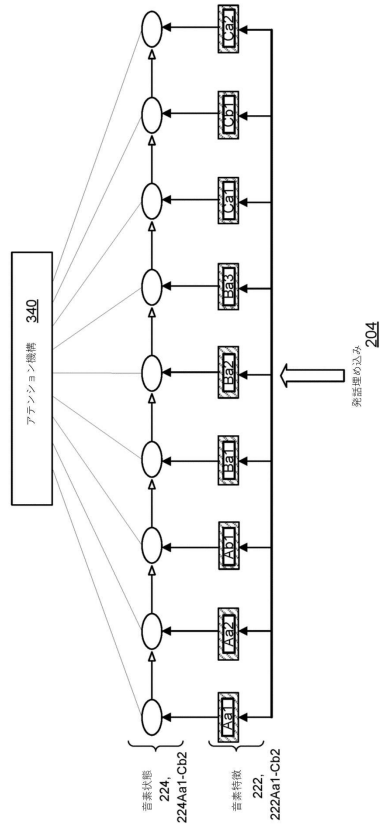


FIG. 3D

10

20

30

40

50

【 図 4 】

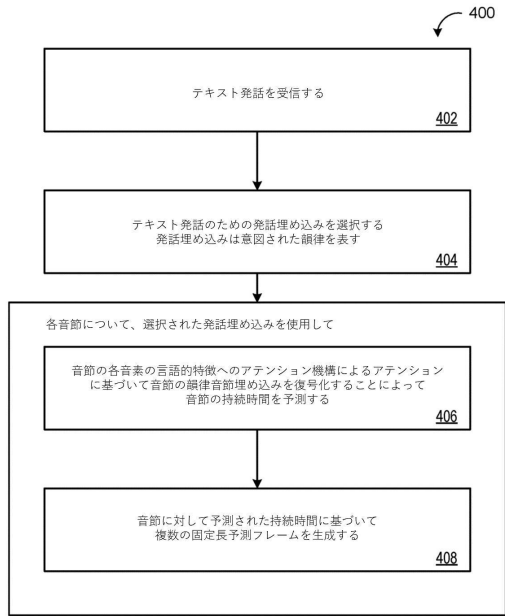


FIG. 4

【 図 5 】

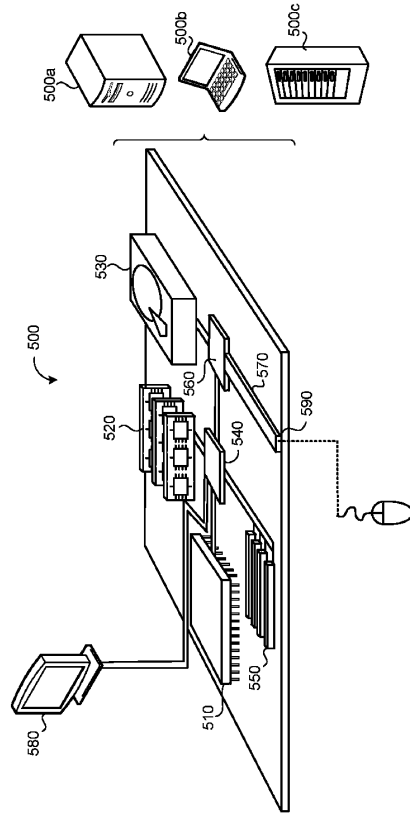


FIG. 5

10

20

30

40

50

---

フロントページの続き

- (72)発明者 チャン、チュン - アン  
アメリカ合衆国 9 4 0 4 3 カリフォルニア州 マウンテン ビュー アンフィシアター パークウ  
エイ 1 6 0 0
- (72)発明者 ワン、ヴィンセント  
アメリカ合衆国 9 4 0 4 3 カリフォルニア州 マウンテン ビュー アンフィシアター パークウ  
エイ 1 6 0 0
- 審査官 山下 剛史
- (56)参考文献 国際公開第 2 0 1 7 / 1 6 8 8 7 0 ( W O , A 1 )  
国際公開第 2 0 1 9 / 2 1 7 0 3 5 ( W O , A 1 )  
田中宏他, V A E - S P A C E : 音声 F 0 パターンの深層生成モデル, 日本音響学会 2 0  
1 8 年春季研究発表会講演論文集 [ C D - R O M ] , 2018年03月, pp.229-230
- (58)調査した分野 (Int.Cl. , D B 名)  
G 1 0 L 1 3 / 0 0 - 1 3 / 1 0 , 2 5 / 3 0  
I E E E X p l o r e